

【特許請求の範囲】

【請求項1】 冗長データ記憶システムにおいてパリティ演算を実行するディスクコントローラであって、実行すべき前記パリティ演算についての情報をそれぞれ識別する複数のタスク記述ブロックを格納するメモリと、前記メモリに接続され、前記複数のタスク記述ブロックのそれぞれにアクセスし、前記タスク記述ブロックで識別される情報に基づいて前記パリティ演算を実行するパリティ演算ロジックと、を備えるディスクコントローラ。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、実行する演算に関連する情報の維持に関し、特にタスク記述ブロックを使用して演算に関連する情報を維持することに関する。

【0002】

【従来の技術】最近の大容量データ格納システムは、データの冗長記憶のために複数の物理ディスクドライブを利用することが多い。このような構成により、データアクセス速度が向上すると共に、任意の単一ディスクの障害により発生しうるデータ損失の保護が容易になる。

【0003】冗長データの記憶には2つの一般的な方法がある。第1の方法、すなわち「ミラーリング」法では、データは複製されて記憶システムの2つの別個の領域に格納される。ディスクアレイには、例えば、同一のデータが2つの別個のディスクに格納される。この方法は、性能が高くまたデータ可用性も高いという利点がある。しかし、ミラーリング法は比較的高価でもあり、データ格納のコストが事実上2倍になる。

【0004】第2の方法、すなわち「パリティ」法では、記憶領域の一部を使用して冗長データを格納するが、冗長記憶領域のサイズは、元のデータの格納に使用される残りの記憶空間よりも小さい。例えば、6枚のディスクを有するディスクアレイでは、5枚のディスクをデータの格納に使用し、6枚目のディスクを「パリティ」データと呼ばれる冗長データの格納専用とする。パリティデータを生存ディスクからのデータと併せて使用することで、1枚のデータディスクからデータを再構築できるようになる。パリティ法は、ミラーリング法よりも費用が低いため有利であるが、同時にミラーリング法と比較して性能が低く可用性も低いという特徴がある。

【0005】パリティ記憶を利用する従来のディスクアレイでは、記憶ディスクにおける空間は、複数の記憶ディスクにわたる複数の記憶ストライプに構成される。各ストライプは複数の記憶空間セグメントからなり、各セグメントは、ディスクアレイの単一記憶ディスク上にあるストライプの一部分である。

【0006】図1は、6つの記憶ディスク13を有する従来のディスクアレイ12を示す。この単純化した例では、

複数の記憶ディスクにわたる5つの記憶ストライプがある。図1は、これら5つのストライプのうちの1つのストライプのデータセグメント及び記憶セグメントを強調して示している。強調して示されたストライプのデータセグメントは、網掛けで示されている。このストライプに対応するパリティセグメントは、黒塗りで示されている。一般に、任意の所与のストライプを構成する6つのセグメントのうち、5つのセグメントはデータセグメントであり、6番目のセグメントがパリティセグメントである。

【0007】この種のパリティ記憶は5+1パリティ記憶と呼ばれ、単一のパリティセグメントにつき5つのデータセグメントがあることを示している。この方式はより一般的にN+1グループ化と呼ばれる。ここで、Nはデータストライプにおけるデータセグメントの実際の数である。

【0008】図1に示すようなN+1冗長グループ化は、任意の1つの物理記憶装置をデータ損失から保護する。記憶装置に障害が発生した場合、生存データからそのデータを再構築することができる。データを回復するために実行される計算は単純であり、この計算は既知である。一般に、次式に従って、データセグメント $D_0 \sim D_{N-1}$ から単一のパリティセグメントPが計算される。

【0009】

【数1】

$$P = x_0 + x_1 + x_2 + \dots + x_{N-1}$$

式中、 $x_0 \sim x_{N-1}$ は、データセグメント $D_0 \sim D_{N-1}$ からのデータに相当する。任意の単一データセグメントを損失した後、同じ式を簡単に変形することでそのデータを回復することができる。

【0010】しかし、多くのシステムでは、2つ以上の記憶装置を損失から保護することが重要になりつつある。そのため、冗長格納システムにおいてN+2グループ化を実施することが必要になりつつある。

【0011】N+2冗長グループ化によりデータ保護が強化されるが、その一方、最初のパリティセグメントの計算及び任意の損失データセグメントの再構築の両方で、より複雑な計算が必要となる。

【0012】N+2パリティ計算の一般式は、次のようになる。

【0013】

【数2】

$$P = p_0x_0 + p_1x_1 + p_2x_2 + \dots + p_{N-1}x_{N-1}$$

$$Q = q_0x_0 + q_1x_1 + q_2x_2 + \dots + q_{N-1}x_{N-1}$$

式中、Pは第1のパリティセグメントの値であり、Qは第2のパリティセグメントの値であり、 $x_0 \sim x_{N-1}$ はデータセグメントの値であり、 $p_0 \sim p_{N-1}$ 及び $q_0 \sim q_{N-1}$ は所与のパリティ方式に固有の定数係数で

ある。

【0014】これらの式は、線形代数の規則により、任意の2つの障害が発生した記憶装置の単一ストライプからのデータを表す任意の2つの未知の $x_a \sim x_b$ について潜在的に解くことができる2つの等式システムを形成する。1つの要件は、2つの係数セット p_i 及び q_i が一次独立していることである。この要件は、例えば、 $p_0 = 1$ 、 $p_1 = 1$ 、 $p_2 = 1$ 、かつ $q_0 = 1$ 、 $q_1 = 2$ 、 $q_2 = 3$ の場合に満たされる。また他の例も可能である。

【0015】 $N+2$ パリティの数学的処理は既知であり、本説明の主題ではない。しかし、上記の簡潔な説明から、 $N+2$ パリティの計算は $N+1$ パリティの計算よりもかなり複雑であることは明白である。 $N+2$ ディスクアレイの実際の実施において、この複雑さは、記憶装置コントローラ及び結果的にディスクアレイ全体のデータスループットを制限する恐れがある。

【0016】

【発明が解決しようとする課題】発明者Michael B. Jacobsonによる「Efficient Parity Operations (効率的なパリティ演算)」という名称の同時係属中の米国特許出願(代理人整理番号10971442)には、 $N+2$ パリティ計算により複雑性が追加されるにも関わらず、適切なデータスループットを維持するより効率的な方法が記載されており、この出願は参照により本明細書に援用される。しかし、この同時係属中の出願に記載のより効率的な方法及び手段を用いても、記憶装置コントローラのデータスループットは、パリティ計算の実行に必要な情報を格納するためのメモリ構造が非効率的であることによりなお悪影響を受けることがある。本発明は、タスク記述ブロックを用いて、パリティ(または他の)演算に関する情報を維持し、適切なデータスループットを維持する方法及び装置を提供する。

【0017】

【課題を解決するための手段】本明細書では、タスク記述ブロックを用いた演算に関連する情報の維持を説明する。

【0018】本発明の一態様によると、タスク記述ブロック待ち行列は複数のタスク記述ブロックそれぞれを指すポインタを格納し、各タスク記述ブロックは1つの演算についての情報を格納する。演算のデータはメモリに格納される。各タスク記述ブロックは、メモリに格納されているデータを指すポインタと、演算ロジックが演算の結果を格納すべきメモリ内の結果位置を指すポインタとを含めた、演算ロジックが単一の演算を実行するために必要な情報を格納する。この情報がすべてタスク記述ブロックに格納されると、タスク記述ブロックを指すポインタがタスク記述ブロック待ち行列に格納され、演算ロジックが要求するときにいつでもその演算を実行することが可能になる。

【0019】パリティ演算に関連する本発明の別の態様によると、タスク記述ブロック待ち行列は複数のタスク記述ブロックそれぞれを指すポインタを格納し、各タスク記述ブロックは1つのパリティ演算についての情報を格納する。パリティ演算についてのデータは、異なるパリティ演算に使用される様々な複数の係数と同様に、メモリに格納される。各タスク記述ブロックは、下記を含め、パリティ演算ロジックが単一パリティ演算を実行するために使用する情報を格納する。すなわち、実行する演算のタイプ、タスク記述ブロックのサイズ、オプションであるタグ番号(一貫性チェックのため)、実行する演算のステータス結果のどれを格納するかについての応答待ち行列の指示、メモリに格納されているデータを指すポインタ、複数の係数のサブセットを指すポインタ、サブセット内の係数の数の指示、及びパリティ演算ロジックがパリティ演算の結果を格納するメモリ内の結果位置を指すポインタを格納する。この情報がすべてタスク記述ブロックに格納されると、タスク記述ブロックを指すポインタがタスク記述ブロック待ち行列に格納され、パリティ演算ロジックが要求するときはいつでもパリティ演算を実行することが可能になる。

【0020】

【発明の実施の形態】<パリティ演算>図2を参照すると、本発明による冗長データ格納システム20は、データストライプ24を有する格納ディスク22を利用する。各データストライプ24は、複数のデータセグメント $x_0 \sim x_{N-1}$ と、少なくとも2つの対応するパリティセグメントP及びQとを含む。P及びQは、データセグメント $x_0 \sim x_{N-1}$ 、第1のパリティ係数セット $p_0 \sim p_{N-1}$ 、及び第2のパリティ係数セット $q_0 \sim q_{N-1}$ から導出される。パリティ係数は以下の式に従い、各データセグメントに対応する。

【0021】

【数3】

$$P = p_0x_0 + p_1x_1 + p_2x_2 + \dots + p_{N-1}x_{N-1}$$

$$Q = q_0x_0 + q_1x_1 + q_2x_2 + \dots + q_{N-1}x_{N-1}$$

【0022】本発明において、パリティ演算は、パリティセグメント生成演算、パリティセグメント再生成演算、及びデータセグメント再構築演算に分類することができる。

【0023】パリティセグメント生成演算は、新しいデータストライプを作成する場合、すなわちパリティセグメントが全く新しいデータに基づいて作成される場合に行われる。

【0024】パリティセグメント再生成演算は、新しいデータにより新しいデータセグメントが追加される場合、または読み出し/変更/書き込みサイクルにより1つまたは複数の既存のデータセグメントが変更される場合に、既存のストライプに対して行われる。パリティセ

グメント再生成演算では、パリティセグメントが、データストライプ全体を読み出すことなく増分的に変更される。例えば、新しいデータにより新しいデータセグメント x_4 が追加されるものと想定する。 P_{NEW} は、次のように計算される。

【0025】

【数4】

$$P_{NEW} = P_{OLD} + p_4 x_4$$

【0026】同様に、読み出し／変更／書き込みサイクルの結果、データセグメント x_2 が変更されるものと想定する。この場合、 P_{NEW} は、次のように計算される。

【0027】

【数5】

$$P_{NEW} = P_{OLD} - p_2 x_{2OLD} + p_2 x_{2NEW}$$

【0028】古いデータセグメント値から新しいP及びQの値を計算するために要するメモリは、ストライプが

$$x_a = f(p_0, p_a) x_0 + f(p_1, p_a) x_1 + \dots + f(p_a) P + \dots + f(p_{N-1}, p_a) x_{N-1}$$

$$x_a = f(q_0, q_a) x_0 + f(q_1, q_a) x_1 + \dots + f(q_a) Q + \dots + f(q_{N-1}, q_a) x_{N-1}$$

式中、 $f()$ は、使用しているパリティ生成コードに固有の適切な係数を生成する変換関数である。

【0032】これらの式の実施形態の1つは、以下の通りである。

【0033】

【数7】

$$x_a = p_a^{-1} (p_0 x_0 + p_1 x_1 + \dots + P + \dots + p_{N-1} x_{N-1})$$

$$x_a = p_a^{-1} (q_0 x_0 + q_1 x_1 + \dots + Q + \dots + q_{N-1} x_{N-1})$$

$$x_a = f(p_0, q_0, p_a, p_b, q_a, q_b) x_0 + f(p_1, q_1, p_a, p_b, q_a, q_b) x_1 + \dots + f(p_a, p_b, q_a, q_b) P + f(p_k, q_k, p_a, p_b, q_a, q_b) x_k + f(p_{k+1}, q_{k+1}, p_a, p_b, q_a, q_b) x_{k+1} + \dots + f(p_a, p_b, q_a, q_b) Q + \dots + f(p_{N-1}, q_{N-1}, p_a, p_b, q_a, q_b) x_{N-1}$$

$$x_b = f(p_0, q_0, p_a, p_b, q_a, q_b) x_0 + f(p_1, q_1, p_a, p_b, q_a, q_b) x_1 + f(p_a, p_b, q_a, q_b) P + f(p_k, q_k, p_a, p_b, q_a, q_b) x_k + f(p_{k+1}, q_{k+1}, p_a, p_b, q_a, q_b) x_{k+1} + \dots + f(p_a, p_b, q_a, q_b) Q + \dots + f(p_{N-1}, q_{N-1}, p_a, p_b, q_a, q_b) x_{N-1}$$

式中、 $f()$ は、使用しているパリティ生成コードに固有の適切な係数を生成する変換関数である。

【0036】これらの式の実施形態の1つは、以下の通りである。

変更される都度初めからP及びQを計算するために要するメモリよりもはるかに少ない。

【0029】本発明によれば、パリティセグメント再生成演算は、データセグメントの追加から生じるパリティ再生成、またはデータセグメントの変更から生じるパリティ再生成のいずれかにさらに分類することができる。

【0030】データセグメント再構築演算は、2つの下位分類、すなわち、単一データセグメント再構築演算と、二重データセグメント再構築演算とを含む。生存データセグメントと組み合わせて、PパリティセグメントまたはQパリティセグメントのいずれかから、単一データセグメントを再構築することができる。一般に、データセグメント x_a は、次のようにしてパリティセグメントPまたはQのいずれかから再構築される。

【0031】

【数6】

【0034】生存データセグメントと組み合わせて、Pパリティセグメント及びQパリティセグメントから、2つのデータセグメントを再構築することができる。一般に、2つのデータセグメント x_a 及び x_b は、次のようにしてパリティセグメントP及びQから再構築される。

【0035】

【数8】

りである。

【0037】

【数9】

$$\begin{aligned}
x_a &= (p, q_b \vdash p_b q_a)^{-1} ((q_b p_0 + p_b q_0)x_0 + (q_0 p_1 + p_0 q_1)x_1 + \dots + q_b P + \dots \\
&\quad + (q_b p_k + p_b q_k)x_k + (q_b p_{k+1} + p_b q_{k+1})x_{k+1} + \dots + p_b Q + \dots + \\
&\quad (q_b p_{N-1} + p_b q_{N-1})x_{N-1}) \\
x_b &= (p_a q_b + p_b q_a)^{-1} ((q_a p_0 + p_a q_0)x_0 + (q_a p_1 + p_a q_1)x_1 + \dots + q_a P + \dots \\
&\quad + (q_a p_k + p_a q_k)x_k + (q_a p_{k+1} + p_a q_{k+1})x_{k+1} + \dots + p_a Q + \dots + \\
&\quad (q_a p_{N-1} + p_a q_{N-1})x_{N-1})
\end{aligned}$$

【0038】一般に、上記パリティ演算はすべて、有限個のかかる係数を有するベースセットから選択される既知の係数の種々の組み合わせを用いることで達成することができる。これら係数は、 $p_0 \sim p_{N-1}$ 、 $q_0 \sim q_{N-1}$ 、及び変換関数 $f(\cdot)$ から得られる係数を含む。任意の特定のパリティ演算は、計算中の実際のデータセグメントまたはパリティセグメントに応じて、これら係数のサブセットを利用する。特定の計算に必要な特定の係数サブセットは、演算の分類及び関係する特定のデータセグメントやパリティセグメントの双方に依存する。従って、所与のパリティ演算の分類内には種々の状況すなわちシナリオがあり、それぞれに異なる係数サブセットが必要となる。例えば、データセグメント x_5 をストライプに追加する場合、係数 p_5 及び q_5 が必要とされる場合に1つのシナリオが発生する。データセグメント x_6 をストライプに追加する場合、係数 p_6 及び q_6 が必要な場合に別のシナリオが発生する。

【0039】＜係数サブセット＞図3は、複数の係数サブセット31を含むメモリアレイ30を示している。各係数サブセットは、対応するデータセグメントに適用されパリティ計算結果を生成する、予め選択された、かつ／または予め計算された係数のリストまたは連結である。本発明によれば、異なる係数サブセットが予め選択され、各異なる演算シナリオについて格納される。次に、パリティ演算ロジックによる参照及び直接の使用のため、サブセットはフォーマットされ線形メモリアレイにパックされる。異なるシナリオには異なる数の係数が必要であるため、サブセットは同じ長さまたはサイズである必要はない。

【0040】本発明の実施形態において、各係数は単一バイトである。「パックされる」という語は、好ましくは使用していない空間に干渉することなく、係数のサブセットまたはストリングが線形メモリ内で連結され、記憶空間を節約することを意味する。

【0041】ここで、パリティ演算を行う場合、サブセット内の係数とデータのセグメント（データセグメントまたはパリティセグメントのいずれか）との間には1対1の対応がある。各係数は、対応するデータセグメントまたはパリティセグメントにのみ適用されて、演算結果を生成する。

【0042】可能な各パリティ計算のケースすなわちシナリオにつき、1つの係数サブセットがアレイに含まれる。特定の各計算シナリオに関するアレイ内のサブセットの開始位置を見つけるために、固有のインデックス式を用いることができる。一般に、サブセットは、P及びQそれぞれに関わる計算に対応する対で構成される。

【0043】図3を参照すると、メモリアレイ30は、複数の分類グループ32、33、34、35、及び36を含み、各分類グループは、特定のパリティ演算分類に対応する係数サブセット31を含む。分類グループにおける各サブセットは、グループの分類内で発生する特定のシナリオに関する係数を有する。1つの例外を除き、係数サブセットは任意の所与の分類グループ内で同一のサイズである。

【0044】アレイの開始位置からグループの開始位置までのグループオフセットを計算することで、アレイ30内で特定の分類グループを見つけることができる。このグループ・オフセットは、グループのアレイへのベースインデックスである。分類グループ内で特定の係数サブセットを見つけるために、グループの開始位置からのサブセット・オフセットがベースインデックスに加えられる。これにより、所望の係数サブセットの開始位置を見つけるアレイ内へのインデックスが生成される。

【0045】本発明の一実施形態によれば、一般的なパリティ演算分類は以下のように定義される。

【0046】1) パリティ生成演算—既存のデータまたはパリティを持たない部分ストライプまたは完全に新しいストライプ。

【0047】2) セグメントの追加によって生じるパリティ再生成演算—新しいデータセグメントを2つのパリティセグメントに組み込むことによる、ストライプの増分的成長。

【0048】3) セグメントの変更によって生じるパリティ再生成演算—すでに2つのパリティセグメントに組み込まれているデータセグメントの変更（読み出し／変更／書き込み）。

【0049】4) 単一データセグメント再構築—2つのパリティセグメントの一方とストライプからの生存データセグメントとを使用した単一データセグメントの再構築。2つの記憶装置に障害が発生する場合に、障害が発生した2つの記憶装置の一方がPまたはQを保持するこ

とができるため、PパリティセグメントまたはQパリティセグメントの何れかからの再構築がサポートされる。

【0050】5) 二重データセグメント再構築—2つのパリティセグメントP及びQと、ストライプからの生存データセグメントとを用いたストライプの2つのデータセグメントの再構築。

【0051】＜分類1の係数サブセットの構造＞アレイの第1の分類グループ32は、パリティ生成演算のための係数サブセットを含む。パリティ生成演算は、新しいPセグメント及びQセグメントを新しいデータセグメント $x_0 \sim x_{N-1}$ から生成する。この分類グループには、2つの係数サブセットしかない。サブセットはそれぞれ、パリティセグメントP及びQの生成に対応している。

【0052】

【数10】

$$P: \{p_0, p_1, \dots, p_{N-1}\} \text{ and}$$

$$Q: \{q_0, q_1, \dots, q_{N-1}\}$$

これらサブセットはそれぞれ同じ長さ (N) である。

【0053】＜分類2の係数サブセットの構造＞アレイの第2の分類グループ33は、ストライプを増分的に追加するパリティ演算のための係数サブセットを含む。このタイプの演算は、任意の所与の新しいまたは追加のデータセグメント $x_a \sim x_b$ (但し、 $b < N$ 、かつ $a \leq b$) の連続範囲と組み合わせてP及びQの各セグメントを更新する。データセグメント $0 \sim N-1$ 内のデータセグメントのあらゆる可能な範囲 $a \sim b$ に対応する、これら演算の異なるシナリオが多数ある。各シナリオには、異なる係数サブセットが必要である。例えば、新しいまたは追加のデータセグメントが x_3 及び x_4 である場合、Pの計算に必要な係数サブセットは $\{p_3, p_4\}$ である。新しいまたは追加のデータセグメントが $x_2 \sim x_5$ である場合、Pの計算に必要な係数サブセットは $\{p_2, p_3, p_4, p_5\}$ である。データセグメント $0 \sim N-1$ 内の可能な範囲の総計は、Nの値に依存する。

【0054】分類グループ2の各係数サブセットは、サブセットがPの計算に適用されるか、またはQの計算に適用されるかを示す2つの初期パラメータを含む。これら初期パラメータはそれぞれ「0」または「1」のいずれかに設定される。これら係数のうちの最初の値「1」は、計算がパリティセグメントPに関わることを示す。これら係数の二番目の値「1」は、計算がパリティセグメントQに関わることを示す。これら2つのパラメータの一方のみが随時「1」に等しく設定されるべきである。

【0055】分類2サブセット内の残りの係数は、新しく追加されたデータストライプからP及びQを再生成するために使用される係数の部分範囲 (sub-range) であ

る。従って、この分類グループは以下の形態の複数の係数サブセットを含む。

【0056】

【数11】

$$P: \{1, 0, p_a, \dots, p_b\} \text{ and}$$

$$Q: \{0, 1, q_a, \dots, q_b\}$$

【0057】分類グループ33は、Nに応じて、より広い範囲の $0 \sim N-1$ 内のあらゆる範囲 $a \sim b$ に対応するような複数のサブセットを含む。アレイのこのセクション内の係数サブセットの長さすなわちサイズは可変であり、各演算シナリオについて $(b-a)$ に等しい。

【0058】この分類グループ内において、係数サブセットは、長さによって配列されグループ化される。すなわち、最小数の係数を含む係数サブセットは、分類グループの最初の部分に配置される。最大数の係数を含む係数サブセットは、分類グループの最後に配置される。これらの各グループの中で、係数サブセットは、係数サブセットでカバーされる範囲の係数の下付文字に従った順に配列される。従って、 $a=0$ を有するサブセットが最初に配置され、 $a=1$ を有するサブセットが次に配置され、以下同様である。

【0059】＜分類3の係数サブセットの構造＞第3の分類グループ34内の係数サブセットは、単一データセグメントが変更されるとき、P及びQを更新するために使用される。このタイプの演算は、データセグメント x_a が変更される場合に、Pセグメント及びQセグメントを更新する。

【0060】分類2のグループのように、各分類3サブセットの最初の2つのパラメータは、グループの係数がPの計算またはQの計算に適用可能であるか否かを示す。これら係数はそれぞれ「0」または「1」に設定される。これら係数の最初の値が「1」であることは、サブセットの係数がパリティセグメントPに適用されることを示す。これら係数の二番目の値が「1」であることは、サブセットの係数がパリティセグメントQに適用されることを示す。

【0061】各サブセットは、変更中のデータセグメント x_a に対応する単一の残りの係数を含む。

【0062】

【数12】

$$P: \{1, 0, p_a\} \text{ and}$$

$$Q: \{0, 1, q_a\}$$

【0063】第3の分類グループ34は、0から $N-1$ までの a のすべての値に対応するN対のこのようなサブセットを含む。これらサブセットは、 $a=b$ である分類2係数サブセットの特別な場合に対応するため、新しい単一データセグメントをストライプに追加する際に使用可能なことに留意されたい。

【0064】＜分類4の係数サブセットの構造＞第4の分類グループ35内の係数サブセットは、パリティセグメントの一方及び生存データセグメントに基づいて単一データセグメント x_a を再構築するために使用される。係数は、選択された誤り修正コードの数学的処理（ f

（ f ）に従って変換されて再構築演算を実行することを除き、分類1の係数に密接に対応する。

【0065】

【数13】

$$P: \{ f(p_0, p_a), f(p_1, p_a), \dots, f(p_a), \dots, f(p_{N-1}, p_a) \}$$

$$Q: \{ f(q_0, q_a), f(q_1, q_a), \dots, f(q_a), \dots, f(q_{N-1}, q_a) \}$$

より具体的には、以下の通りである。

【数14】

【0066】

$$P: (p, q_b + p_b q_a)^{-1} ((q_b p_0 + p_b q_0), (q_0 p_1 + p_0 q_1), \dots, q_b P, \dots, (q_b p_k + p_b q_k), (q_b p_{k+1} + p_b q_{k+1}), \dots, p_b Q + \dots + (q_b p_{N-1} + p_b q_{N-1}))$$

$$P: (p_a q_b + p_b q_a)^{-1} ((q_a p_0 + p_a q_0), (q_a p_1 + p_a q_1), \dots, q_a P + \dots + (q_a p_k + p_a q_k), (q_a p_{k+1} + p_a q_{k+1}), \dots, p_a Q, \dots, (q_a p_{N-1} + p_a q_{N-1}))$$

【0067】第4の分類グループは、0から $N-1$ までの a のすべての値に対応する N 対のこのようなサブセットを含む。各サブセットにおいて、係数 $f(p_a)$ または係数 $f(q_a)$ はデータセグメント x_a に対応することに留意されたい。

セグメント及び生存データセグメントに基づいて、2つのデータセグメント x_a 及び x_b を再構築するために使用される。係数は、選択された誤り修正コードの数学的処理（ f （ f ））に従って変換されて再構築演算を実行することを除き、分類1の係数に密接に対応する。

【0069】

【数15】

【0068】＜分類5の係数サブセットの構造＞第5の分類グループ36内の係数サブセットは、2つのパリティ

$$x_a: \{ f(p_0, q_0, p_a, p_b, q_a, q_b), f(p_1, q_1, p_a, p_b, q_a, q_b), \dots, f(p_a, p_b, q_a, q_b), \dots, f(p_k, q_k, p_a, p_b, q_a, q_b), f(p_{k+1}, q_{k+1}, p_a, p_b, q_a, q_b), \dots, f(p_a, p_0, q_a, q_b), \dots, f(p_{N-1}, q_{N-1}, p_a, p_b, q_a, q_b) \}$$

$$x_b: \{ f(p_0, q_0, p_a, p_b, q_a, q_b), f(p_1, q_1, p_a, p_b, q_a, q_b), \dots, f(p_a, p_b, q_a, q_b), \dots, f(p_k, q_k, p_a, p_b, q_a, q_b), f(p_{k+1}, q_{k+1}, p_a, p_b, q_a, q_b), \dots, f(p_a, p_0, q_a, q_b), \dots, f(p_{N-1}, q_{N-1}, p_a, p_b, q_a, q_b) \}$$

【0070】アレイの5番目のセクションは、0から $N-1$ の範囲内の a 及び b のあらゆる可能な組み合わせに対応する、 $(N \times (N-1)) / 2$ 対のかかるサブセットを含む。各サブセットにおいて、係数 $f(p_a, p_b, q_a, q_b)$ は、データセグメントのうちどれが再構築されているかに応じて、データセグメント x_a ま

たは x_b に対応することに留意されたい。

【0071】これらの式の可能な実施形態の1つは、以下の通りである。

【0072】

【数16】

$$x_a: (p_a, q_b + p_b q_a)^{-1} (q_b, p_0 + p_b q_0), (p_a, q_b + p_b q_a)^{-1} (q_b, p_1 + p_b q_1), \dots, (p_a, q_b + p_b q_a)^{-1} q_b, \dots, (p_a, q_b + p_b q_a)^{-1} (q_0, p_k + p_0 q_k), (p_a, q_b + p_b q_a)^{-1} (q_0, p_{k+1} + p_0 q_{k+1}), \dots, (p_a, q_b + p_b q_a)^{-1} p_b, \dots, (p_a, q_b + p_b q_a)^{-1} (q_0, p_{N-1} + p_0 q_{N-1})$$

$$x_b: (p_a, q_b + p_b q_a)^{-1} (q_b, p_0 + p_b q_0), (p_a, q_b + p_b q_a)^{-1} (q_a p_1 + p_a q_1), \dots, (p_a, q_b + p_b q_a)^{-1} q_a, \dots, (p_a, q_b + p_b q_a)^{-1} (q_a, p_k + p_a q_k), (p_a, q_b + p_b q_a)^{-1} (q_a, p_{k+1} + p_a q_{k+1}), \dots, (p_a, q_b + p_b q_a)^{-1} p_a, \dots, (p_a, q_b + p_b q_a)^{-1} (q_a, p_{N-1} + p_a q_{N-1})$$

【0073】＜係数サブセットの使用＞図4は、上記アレイ格納方式に従ってパリティ演算を実行する方法を示す。第1のステップ100で、種々のパリティ演算をパリティセグメント生成演算、パリティセグメント再生成演算、及びデータセグメント再構築演算を含む分類に分類する。より具体的には、演算は、パリティ生成演算、セグメント追加によって生じるパリティ再生成演算、セグメントの変更によって生じるパリティ再生成演算、単一データセグメント再構築演算、または二重データセグメント再構築演算のいずれかに分類される。パリティ演算の各分類は、複数の多様な分類シナリオを含み、その各シナリオは各パリティ係数サブセットに関係する。

【0074】ステップ102で、個々のパリティ係数を予め計算し、異なるパリティ演算及びパリティ演算の異なるシナリオに使用するパリティ係数サブセットを予め選択する。このステップは、すでに述べた説明に従って行われる。

【0075】ステップ104で、予め選択したパリティ係数サブセットのすべてを、インデックスの付いた線形メモリアレイに格納する。サブセットには、パリティ計算ロジックによりアクセス可能である。このステップは、係数サブセットを予めフォーマットして、ハードウェアベースのパリティ演算ロジックが効率的に利用できるようにすることを含む。特に、各サブセットの個々の係数は隣接するバイトまたは記憶装置にパックされ、ハードウェアベースのパリティ演算ロジックに固有の様式に配列される。このステップの結果、メモリアレイは、種々の計算シナリオそれぞれに対応する単一の係数サブセットを含むようになる。

【0076】個々の係数及び係数サブセットは、非干渉データ要素を用いてパックされる。アレイのサブセットはグループ化され、すでに述べたように順序づけられ、係数サブセットはそれらの分類順によって分類グループにグループ化される。第2の分類グループ内では、サブセットは可変サイズを有する。さらに、第2の分類グループ内のサブセットは、サイズによりサブグループ化され、最小の番号を有する係数に従い昇順に順序づけられ

$$(((L_i - 1)(12N + L_i(3N - 2L_i - 5))/6) - 3(N - 1)) + a(L_i + 2)$$

第2の分類グループのサイズは、次の式によって与えられる。

【0082】

【数19】

$$((N - 1)(12N + N(3N - 2N - 5))/6) - 3(N - 1)$$

【0083】メモリアレイへの適切なオフセットを決定した後、決定したパリティ係数サブセットをメモリから読み出すステップ108が行われる。ステップ110で、メモリから読み出されたパリティ係数サブセットを用いて、特定のパリティ演算を実行する。

【0084】＜ディスクコントローラの動作＞図5は、本発明によるディスクコントローラ200の最も適切な構

成要素を示す。

【0077】パリティ演算中、パリティ演算ロジックがメモリアレイにアクセスし、パリティ演算の種々のシナリオでの使用に適した係数サブセットを得る。よって、ステップ106で、格納されているパリティ係数サブセットのうちどれが特定のパリティ演算に必要なのかを決定する。このステップは、パリティ演算の分類と線形メモリアレイへのグループオフセットの決定を含み、そのパリティ演算分類に対応する分類グループの開始位置を示す。次に、グループ内への所望の係数サブセットの位置までのサブセットオフセットが計算される。

【0078】ステップ106は、第2の分類グループ以外単独である。詳細に上述したように、第2の分類グループは長さすなわちサイズが可変の係数サブセットを含み、特定の係数サブセットのオフセットを決定することが困難である。しかし、本願発明者らは、第2の分類グループが上述したように配列されるとき、同じサイズの係数サブセットのサブグループを順序づけると、特定のサブグループに対するオフセットを、係数サブセットのサブグループのサイズ及びN（任意のサブグループに含まれる最大数の係数）の関数として計算することができることを見出した。具体的には、サブセットサイズ L_i に対応する特定のサブグループiに対するオフセットは、以下に等しい。

【0079】

【数17】

$$((L_i - 1)(12N + L_i(3N - 2L_i - 5))/6) - 3(N - 1)$$

【0080】この式は、各サブセットには、P及びQに対応する定数の対（上述）が存在するものと想定している。しかし、 L_i は $(b - a)$ に等しい。サブグループi内において、特定の係数サブセットのオフセットは a ($L_i + 2$)に等しい。従って、 $x_a \sim x_b$ に対応する係数の範囲に関して、分類グループへの全体的なオフセットは、次のようになる。

【0081】

【数18】

成要素を示す。ディスクコントローラ200は、マイクロプロセッサ201及び関連するメモリ202を有する。さらに、ディスクコントローラ200は、ハードディスクインタフェースコンポーネント203及び通信コンポーネント204を備える。ハードディスクインタフェースコンポーネント203は、ディスクコントローラ200に接続されて制御されるハードディスクへのアクセス手段を提供する。通信コンポーネント204は、ホストコンピュータとハードディスクコントローラ200の間のインタフェースとして機能する。

【0085】これらコンポーネントに加え、ハードディスクコントローラ200は、特定用途向け集積回路（AS

IC)の形態のハードウェアベースのパリティ演算ロジック205を備える。「ハードウェアベース」という語は、このロジックコンポーネントが、ソフトウェアベースのロジックとは対照的に、プログラムメモリからの命令の検索や実行を行わないことを意味する。むしろ、ロジックコンポーネントは、信号及びデータを処理する専用の相互接続されるロジック要素を有する。このようなハードウェアベースのロジックはマイクロプロセッサまたは他の命令ベースのプロセッサよりも柔軟性が低い。が、ハードウェアベースのロジックは命令ベースのロジックよりもはるかに高速であることが多い。

【0086】一般に、ディスクコントローラは以下のように動作する。マイクロプロセッサ201は、ホストコンピュータとの通信を処理し、ホストコントローラへの、かつホストコントローラからのすべてのデータ転送を調整する。さらに、マイクロプロセッサ201は、実際のディスク転送をすべて調整する。しかし、データは、ディスクに書き込まれる前にメモリ202にバッファリングされる。パリティ演算は、マイクロプロセッサ201の制御下でメモリ202内のデータに対して行われる。

【0087】初期化中、マイクロプロセッサ201は、メモリ202内に係数サブセットテーブル212を構築する。続いて、パリティ演算のときになると、マイクロプロセッサ201は、パリティ演算の分類及びシナリオを決定する。この情報がいったん決定されると、マイクロプロセッサ201は、パリティ演算の目的となる1つまたは複数のデータセグメント及びパリティセグメント（まとめてデータブロック214と呼ぶ）の、メモリ202における位置を示すスクリプトを作成する。スクリプトは、パリティ演算に適切な係数サブセットが見つかる係数サブセットテーブル212へのオフセット、及び係数サブセットに含まれる係数の数を示す。スクリプトはまた、要求された計算の結果が配置されるメモリ202内の位置も示す。各スクリプトは、単一パリティ演算についての情報を格納し、かかるスクリプトを格納するためのメモリ構造は、本明細書においてタスク記述ブロック（TDB）と呼ばれる。TDBはメモリ202内の特定の位置に格納され、その位置（例えば、32ビットアドレス）へのポインタはメモリ202内のTDB待ち行列216に格納される。応答待ち行列218もまたメモリ202に格納され、ハードウェアロジック205がこれを用いて、パリティ演算が成功したか否かについての指示をマイクロプロセッサ201に戻す。

【0088】スクリプトがメモリに配置されると、待ち行列216におけるスクリプトに対するTDBへのポインタが存在することで、ハードウェアロジックに通知される。ハードウェアロジックはこれにตอบสนองして、(a)指示された係数、データセグメント、及びパリティセグメントを検索し、(b)指示された係数に基づいて適切なパリティ演算を行い、(c)データ及び／または計算し

たパリティセグメントをメモリに戻し、(d)エントリを応答待ち行列218に書き込むことで、演算が首尾良く完了したか否かをマイクロプロセッサ201に知らせる。

【0089】ハードウェアロジックは、係数及びデータ／パリティセグメントの積を合計することで、各種の異なるパリティ演算を行うように構成される。異なる演算は実際に、係数、データセグメント、及びパリティセグメントの数及び選択においてのみ異なる。これら変数は、スクリプトによって特定される。このように、この演算は、ハードウェアベースの計算に非常に都合が良い。

【0090】図6は、タスク記述ブロック250の例示的な構造をさらに詳細に示す。本明細書において、TDB 250は、TDB 250内の情報に基づいてパリティ演算を行うために、ハードウェアロジック205によりアクセスされることに関連して説明される。代わりに、図6に示すメモリ構造は、種々の方法で（例えば、ソフトウェアまたはファームウェアを実行するマイクロプロセッサまたは他の命令ベースのプロセッサにより）パリティ演算を実行する他のシステムで使用することもできる。

【0091】TDB待ち行列216は複数のエントリを含み、各待ち行列エントリは、要求ヘッダ及びTDBポインタの両方を含む。TDBポインタ252はTDB 250をポイントする（例えば、TDB 250の開始アドレスをメモリ202に格納する）、対応する要求ヘッダ251は、TDB 250に基づいて実行される要求された動作についての様々な情報を格納する。一実施形態においては、TDBポインタ252及び要求ヘッダ251の両方が32ビット値であり、要求ヘッダ251は以下の情報を含む。すなわち、TDB 20のサイズ（例えば、8バイトダブルワードグループで）、ハードウェアロジック205がTDB 250内のデータに基づいて実行すべき演算の指示、タグベースの一貫性チェック（より詳細に後述する）を実行すべきか否か、演算完了情報を書き込むべき応答待ち行列218の識別、及びタグ番号（タグベースの一貫性チェックが実行中である場合）である。

【0092】TDB 250は、ヘッダ253と、複数のディスクデータポインタ254と、出力ポインタ256と、係数ポインタ及びデータ258とを含む複数のフィールドを有する。ヘッダ253は、TDBについての様々な情報を含み、一実施形態では以下を含む。すなわち、TDB 250のサイズ（例えば、8バイトダブルワードグループで）、ハードウェアロジック205がTDB 250内のデータに基づいて実行すべき演算の指示（例えば、上記パリティ演算の5つの分類のうちの1つの識別）、ハードウェアロジック205がTDB 250内のデータに基づいて実行すべき演算の実行優先度（例えば、所望であれば、異なるタイプの演算に異なる実行優先度を割り当てることができる）、及びタグ番号（タグベースの一貫性チェックが

実行中である場合)である。

【0093】各ディスクデータポイント254は、パリティ演算に関するディスクデータのメモリ202内の位置をポイントする(図5のデータブロック214のうちの1つをそれぞれポイントする)。実行中のパリティ演算に応じて、各ポイント254によってポイントされるデータは、パリティ演算を行うべきストライプのデータセグメントまたはパリティセグメントである。ポイント254は、ハードウェアロジック205により(ストライプにおいて最も小さい番号の付いたデータセグメントから最も大きな番号の付いたデータを配列し、その後Pパリティセグメント及びQパリティセグメントが続くといった)既知の順に従って配列されるため、ハードウェアロジック205は、パリティ演算についてどのメモリ位置にどのセグメントデータがあるかを知る。出力ポイント256は、パリティ演算ロジック205がパリティ演算の結果を書き込むべきメモリ202内の位置260を指すポイントである。実行中のパリティ演算のタイプに応じて、演算の結果はPパリティ値及びQパリティ値のいずれか(例えば、パリティ生成演算またはパリティ再生成演算の場合)、またはデータセグメント値(例えば、データセグメント再構築の場合)であることができる。単一の出力ポイント256を図6に示しているが、代わりに複数の出力ポイントをTDB250に含めてもよい(例えば、P値に1つのポイント、Q値に1つのポイントというように、複数の値が生成される場合)。

【0094】係数ポイント及びデータ258は、パリティ演算ロジック205が使用すべき係数サブセットが位置する、係数テーブル30内の位置を指すポイントを含む。ポイント及びデータ258内には、係数テーブルから検索すべき係数の数の指示も存在し得る(これは、ヘッダ253において識別される実行すべきパリティ演算のタイプの指示とは別のものであってもよく、また代わりに実行すべきパリティ演算のタイプの指示に埋め込んでも(または固有であっても)よい)。従って、係数ポイント及びデータ258は(ヘッダ253と組み合わせて)、どこでパリティ演算に関する係数を得るか、及びどのくらい係数を得るべきか、ならびにどのパリティ演算を実行すべきかという必要な情報をパリティ演算ロジック205に与える。

【0095】必要な情報がすべてパリティ演算のTDB250に格納されると、マイクロプロセッサ201は、そのTDBを指すポイント及び対応する要求ヘッダ情報を含むエントリをTDB待ち行列216に追加する。待ち行列216に演算を加え、必要なデータが適切な位置で利用できるようになるまで待ち行列216へのエントリの追加を行わないことにより、パリティ演算ロジック205は、確実にパリティ演算の実行に必要な情報を有する。従って、パリティ演算ロジック205は、他のパリティ演算が待ち行列で待機している間、パリティに関する追加情報を待つ

ことで遅延されず、ロジック205の高速演算が可能になる。

【0096】パリティ演算を実行する前に、パリティ演算ロジック205は、一貫性チェックをオプションとして実行して、正しいTDBを使用していることを検証することができる。一実施形態では、この一貫性チェックは、演算のタイプ及びTDBのサイズの検証に基づく。この実施では、TDB250のヘッダ253は、上記と同様に、実行すべき演算のタイプの指示及びTDB250のサイズを含む。正しいTDBにアクセスしていることを検証するために、ハードウェアロジック205は、要求ヘッダ251内の演算のタイプ及びTDBサイズにアクセスし、これらがヘッダ253に格納されている演算のタイプ及びTDBサイズに一致するか否かをチェックする。ハードウェアロジック205が、2つの演算タイプ及び2つのTDBサイズが一致すると判断した場合、TDB250は正しいTDBであると検証され、その他の場合、TDBは正しいTDBではなく、ハードウェアロジック205は、(例えば、応答待ち行列218内のエントリを介して)コマンド不一致のためにパリティ演算が失敗したことをマイクロプロセッサ201に知らせる。

【0097】一貫性チェックはさらに行うことが可能であり、本明細書においてこれはタグベースの一貫性チェックと呼ばれる。タグベースの一貫性チェックを用いる場合、マイクロプロセッサ201はタグ番号を演算に割り当て、タグ番号は要求ヘッダ251及びヘッダ253の両方に格納される。ハードウェアロジック205は、演算タイプ及びTDBサイズを検証することに加え、ヘッダ251及び253におけるタグ番号の一致もチェックする。タグ番号が一致する場合、TDB250は正しいTDBであると検証され(演算タイプ及びTDBサイズも一致するものと仮定する)、その他の場合は、TDBは正しいTDBではなく、ハードウェアロジック205は、(例えば、応答待ち行列218内のエントリを介して)タグ不一致によりパリティ演算が失敗したことをマイクロプロセッサ201に知らせる。

【0098】ハードウェアロジック205が演算を完了すると、計算した演算結果が、出力ポイント256によってポイントされるメモリ位置260に格納される。さらに、ハードウェアロジック205によって実行された演算のステータスを示すエントリが応答待ち行列218に追加される。複数の応答待ち行列218が存在する場合があり、ハードウェアロジック205がエントリを追加する適切な応答待ち行列は、演算に対応する要求ヘッダ251で示される。

【0099】各応答待ち行列218は複数のエントリを含むことができ、各エントリは応答ヘッダ及びTDBポイントを含む。TDBポイント262は、TDB250をポイントするが(例えば、TDB250の開始アドレスをメモリ202に格納する)、対応する応答ヘッダ261は、TDB250

に基づいて実行された演算についての様々な情報を格納する。一実施形態では、TDBポインタ262及び応答ヘッダ261の両方が32ビット値であり、応答ヘッダ261は以下の情報を含む。すなわち、TDB250のサイズ（例えば、8バイトダブルワードグループで）、ハードウェアロジック205がTDB250内のデータに基づいて実行した演算の指示、タグベースの一貫性チェックが実行されたかどうか、二重パリティブロック計算の結果、ゼロしか含まないPパリティブロック及びQパリティブロックが生じたり、または単一パリティブロック計算の結果ゼロしか含まないPパリティブロックが生じる場合に設定される「パリティ・オール・ゼロ」フラグ（このフラグは、演算がパリティ演算ではない場合でも設定される）、ソースデータが定義されたフォーマットパターンに一致する場合に設定される「パターン一致」フラグ、演算を実行する際にどのタイプのエラーが（もしあれば）発生したかに関する指示、及びタグ番号（タグベースの一貫性チェックが実行中である場合）を含む。一実施形態では、応答ヘッダ261を介して以下のエラーを示すことができる。すなわち、タグ不一致（待ち行列216のヘッダ251内のタグが、TDB250のヘッダ253内のタグと一致しない）、コマンド不一致（ヘッダ251における演算タイプ及びTDBサイズが、ヘッダ253内の演算タイプ及びTDBサイズに一致しない）、パリティエラー（メモリ202から戻されるデータにおいて検出されたパリティエラー）、メモリテストデータ不一致（mismatch）（パリティ演算ではなくメモリテスト演算が実行中である場合、これは、メモリテスト中メモリ202から戻されるデータが不正確であったことを示す）、及び不正コマンド（ヘッダ251及びヘッダ253で示される演算タイプが正当な値ではない）である。

【0100】従って、図6に示すメモリ構造を使用して、単一TDBを用いてパリティ演算ロジック205がパリティ演算を実行するために必要な情報をすべて識別する。単一ポインタを用いて各TDBを識別することができ、複数のパリティ演算を図5のTDB待ち行列216に容易に入れることが可能になる。次いで、パリティ演算ロジック205は、演算を待ち行列216から取り除くことにより、パリティ演算を特定の様式で（例えば、先入れ先出し（FIFO）方式で）処理することができる。

【0101】さらに、図5及び図6に示すように、どの演算を実行する必要があるかについての記述は、作業を実際に実行するコンポーネントから分離される。この分離により、2つの異なるプロセス（何をすべきかを記述するプロセス及びする必要があることを実際に実行するプロセス）を、他方がどのように実行されているか（あるいは、どのコンポーネントまたはエンティティがそのプロセスを実行しているか）に関わらず、互いに独立して動作することが可能になる。従って、図5のマイクロプロセッサ201は、どのコンポーネントが実際にパリティ

演算を実行しているのかに関わらず（例えば、パリティ演算は、ハードウェアロジック205または代替として他のコンポーネント（ソフトウェア命令を実行する別のマイクロプロセッサ等）によって実行することができる）、パリティ演算についての必要な情報を準備することができる。同様に、ハードウェアロジック205は、パリティ演算についての情報をメモリ202から獲得でき、どのエンティティが情報をメモリ202内に準備するか、または情報がどのように得られたかに関わらず、パリティ演算を実行することができる。

【0102】図7は、上記メモリ構造によるパリティ演算についてのタスク記述ブロックを使用する方法を示す。図7の方法は、図5のマイクロプロセッサ201及びハードウェアロジック205によって実行される動作を示す（マイクロプロセッサ201によって実行される動作は左側にあり、ハードウェアロジック205によって実行される動作は右側にある）。第1のステップ270で、実行するパリティ演算のTDBを生成する。この生成は、様々な従来の方法の任意のもので実行することができる（例えば、TDBにメモリを割り当てる、既に割り当てられたが使用されていないメモリを使用するなど）。ステップ272で、コマンドタイプ、応答待ち行列識別子、及びタグ番号（オプション）をTDBのヘッダに格納する。他の情報（上述）がヘッダに含まれてもよい。ステップ274で、パリティ演算に関するディスクデータメモリ位置を指すポインタをTDBに格納し、ステップ276で、パリティ演算に関する結果メモリ位置を指すポインタをTDBに格納する。ステップ278で、使用に適した係数サブセットを決定し（上記図4のステップ106と同様）、ステップ280で、その係数サブセットを指すポインタをTDBに格納する。ステップ282で、パリティ演算ロジックにどのパリティ演算を実行すべきかを知らせるために、パリティ演算分類の指示をTDBに格納する。これらの情報すべてがTDBに格納されると、ステップ284で、TDB待ち行列に、要求ヘッダ及びこの新たに追加されたTDBを指すポインタを含むエントリを追加する。従って、TDBの指示は、パリティ演算についてのすべてのデータがメモリ202内で利用できるようになるまで、TDB待ち行列に追加されず、これによって、パリティ演算ロジックがメモリ202にアクセスし、所望のときにいつでもパリティ演算を実行することが可能になる。

【0103】TDB待ち行列エントリが追加された後のある時点において、ステップ286で、ハードウェアロジック205が要求された演算を実行する。演算が完了または失敗すると、ステップ288で、演算が首尾良く完了したか、それともエラーが発生したかを示す（上記ステップ272において識別されるような）応答エントリを適切な応答待ち行列に追加する。

【0104】図7の説明は、TDBエントリ（要求ヘッ

ダ及びTDBポインタの両方)が同時にTDB待ち行列に追加される(ステップ284)ことを述べている。代わりに、情報が利用可能な場合はTDB待ち行列内の要求ヘッダを埋めても良く、この場合は、ステップ284で、TDBポインタだけを待ち行列エントリに追加する(及び本質的にTDB待ち行列エントリが完全に埋められたことをハードウェアロジック205に通知する)。

【0105】本発明について、主に異なるタイプのパリティ演算を冗長ディスクドライブシステムにおいて実行することに関連して説明したが、本発明は、代替的に他の状況で 사용할 ことができる。例えば、本発明は、ハードウェア加速メモリテストを実行するために使用することも可能であり、この場合、検証のために書き込まれるべきメモリアドレスは、図6のフィールド254によって識別され、図5のハードウェアロジック205がテストデータを識別されたメモリアドレスに書き込み、そのメモリアドレスからデータを読み戻し、読み戻されたデータをポインタ256によって識別される位置に格納する。

【0106】別の例として、本発明は、コピー機能を実行するために使用することができる。この場合、コピーするデータはポインタ254によって識別され、得られるコピーの位置はポインタ256によって識別され、ハードウェアロジック205がコピーを実行する。

【0107】さらに別の例において、本発明は、パターンマッチングを実行するために使用することもできる。この場合、比較されるデータセグメントはポインタ254によって識別され、ハードウェアロジック205が比較を実行し、比較の結果(例えば、2つのデータセグメントが一致するかどうかの指示)をポインタ256によって識別される位置に書き込む。

【0108】さらに別の例として、本発明は、チェックサム計算の実行に使用することもできる。この場合、チェックサムが計算されるデータセグメントがポインタ254によって識別され、ハードウェアロジック205がチェックサムを計算し、得られるチェックサム値をポインタ256によって識別される位置に書き込む。

【0109】本発明には例として以下の実施形態が含まれる。

【0110】1. 冗長データ記憶システムにおいてパリティ演算を実行するディスクコントローラ(200)であって、実行すべき前記パリティ演算についての情報をそれぞれ識別する複数のタスク記述ブロック(250)を格納するメモリ(202)と、前記メモリ(202)に接続され、前記複数のタスク記述ブロック(250)のそれぞれにアクセスし、前記タスク記述ブロック(250)で識別される情報に基づいて前記パリティ演算を実行するパリティ演算ロジック(205)と、を備えるディスクコントローラ。

【0111】2. 前記メモリ(202)は、それぞれが前記複数のタスク記述ブロック(250)のうちの1つを識

別する複数のエントリを有するタスク記述ブロック待ち行列(216)をさらに格納し、前記パリティ演算ロジック(205)は、前記複数のエントリのうちの1つによって識別されるタスク記述ブロック(250)にアクセスすることによって、パリティ演算の実行に必要な情報を得る、上記1に記載のディスクコントローラ(200)。

【0112】3. 前記複数のタスク記述ブロック(250)はそれぞれ、対応するパリティ演算を実行するために必要な情報を格納し、前記情報には、前記パリティ演算に関するディスクデータについて、前記メモリ(202)内の位置をそれぞれ識別する複数のディスクデータポインタ(254)と、前記パリティ演算の結果を格納すべき前記メモリ(202)内の位置を識別する出力ポインタ(256)と、前記パリティ演算に関するパリティ係数が格納される前記メモリ(202)内の位置を識別する係数ポインタ(258)と、実行すべき前記パリティ演算のタイプの識別子(253)と、が含まれる、上記1に記載のディスクコントローラ。

【0113】4. パリティ演算についての情報をメモリ(202)に格納するステップ(272、274、276、280、282)であって、前記パリティ演算は、冗長データ記憶システムにおけるストライプについて少なくとも1つのパリティセグメントを生成するか、または前記ストライプについて少なくとも1つのデータセグメントを再生成するものである、ステップと、前記パリティ演算に必要な前記情報がすべて前記メモリ(202)に格納された後でのみ、パリティ演算ロジック(205)が前記情報を利用できるようにするステップと、を含む方法。

【0114】5. 前記格納ステップは、前記パリティ演算についての前記情報を前記メモリ(202)に格納するステップと、前記情報を指すポインタを前記パリティ演算に対応するタスク記述ブロック(250)に格納するステップと、を含む上記4に記載の方法。

【0115】6. 前記格納ステップは、前記情報のメモリ位置識別子を前記パリティ演算に対応するタスク記述ブロック(250)に格納するステップを含み、前記情報を利用できるようにするステップは、前記タスク記述ブロック(250)の識別子を、前記パリティ演算ロジック(205)がアクセス可能なタスク記述ブロック待ち行列(216)に追加するステップを含む、上記4に記載の方法。

【0116】7. 前記格納ステップは、前記メモリ(202)内にタスク記述ブロック(250)を生成するステップ(270)を含み、該タスク記述ブロック(250)には、前記パリティ演算に関するディスクデータの、前記メモリ(202)内の位置をそれぞれ識別する複数のディスクデータポインタ(254)と、前記パリティ演算の結果を格納すべき前記メモリ(202)内の位置を識別する出力ポインタ(256)と、前記パリティ演算に関するパリティ係数が格納される前記メモリ(202)内の位置を識別す

る係数ポインタ (258) と、実行すべき前記パリティ演算のタイプの識別子 (253) と、を含む、上記 4 に記載の方法。

【0117】8. ハードウェア演算ロジック (205) が情報を利用できるようにして、1つまたは複数のデータセグメントに基づいて演算を実行する方法であって、演算に関する1つまたは複数のデータセグメントをそれぞれメモリ (202) 内の別々の位置に格納するステップと、前記演算に対応するタスク記述ブロック (250) に、前記メモリ (202) 内の別々の位置それぞれへのポインタ (254) を格納するステップと、前記タスク記述ブロック (250) に、前記演算の結果を格納すべき前記メモリ (202) 内の結果位置へのポインタ (256) を格納するステップと、タスク記述ブロック待ち行列 (216) に、前記タスク記述ブロック (250) へのポインタ (252) を格納するステップと、を含む方法。

【0118】9. 前記ハードウェア演算ロジック (205) が前記パリティ演算を実行し、前記方法は、前記メモリ (202) に異なる前記パリティ演算に使用される複数の係数を格納するステップと、前記タスク記述ブロック (250) に、前記パリティ演算ロジック (205) が前記パリティ演算に使用する複数の係数サブセットへのポインタ (258) を格納するステップと、をさらに含む、上記 8 に記載の方法。

【0119】10. 前記ポインタ (258) を前記タスク記述ブロック (250) に格納するステップは、前記メモリ (202) 内の前記複数の異なる位置それぞれを指す前記ポインタ (254) が前記タスク記述ブロック (250) に格納された後でのみ、かつ前記結果位置を指す前記ポインタ (256) が前記タスク記述ブロック (250) に格納さ

れた後でのみ、前記タスク記述ブロック (250) を指す前記ポインタを前記タスク記述ブロック待ち行列 (216) に格納するステップを含む、上記 8 に記載の方法。

【図面の簡単な説明】

【図1】従来技術によるN+1冗長グループ化を示すブロック図である。

【図2】本発明の実施形態において使用されるN+2冗長グループ化を示すブロック図である。

【図3】本発明によるメモリテーブルのレイアウトを示すブロック図である。

【図4】本発明の特定の実施形態によりパリティ演算を実行する方法を示す図である。

【図5】本発明によるディスクコントローラの関連コンポーネントを示すブロック図である。

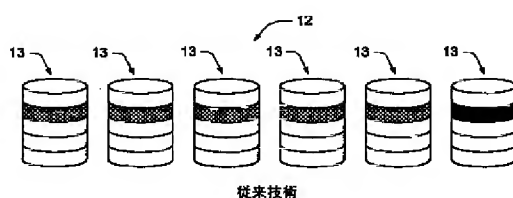
【図6】タスク記述ブロックの例示的な構造をさらに詳細に示す図である。

【図7】本発明の特定の実施形態によるタスク記述ブロックを使用する方法を示す図である。

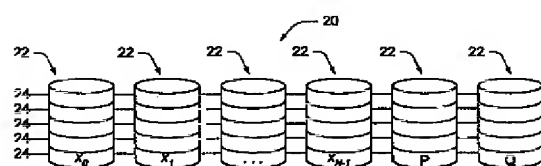
【符号の説明】

- 200 ディスクコントローラ
- 202 メモリ
- 205 演算ロジック
- 216 タスク記述ブロック待ち行列
- 250 タスク記述ブロック
- 252 タスク記述ブロックポインタ
- 253 ヘッダ
- 254 ディスクデータポインタ
- 256 出力ポインタ
- 258 係数ポインタ及びデータ

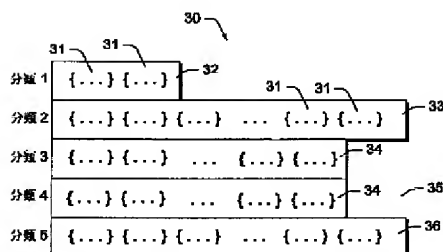
【図1】



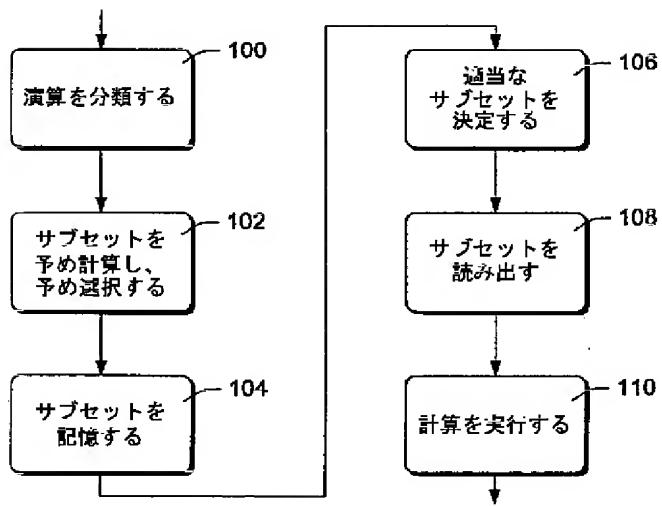
【図2】



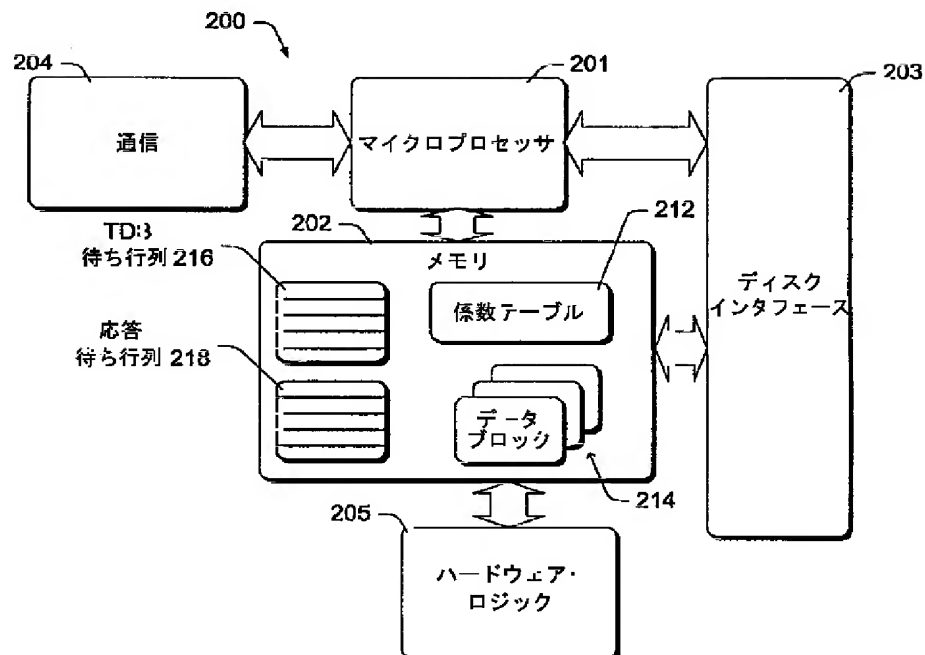
【図3】



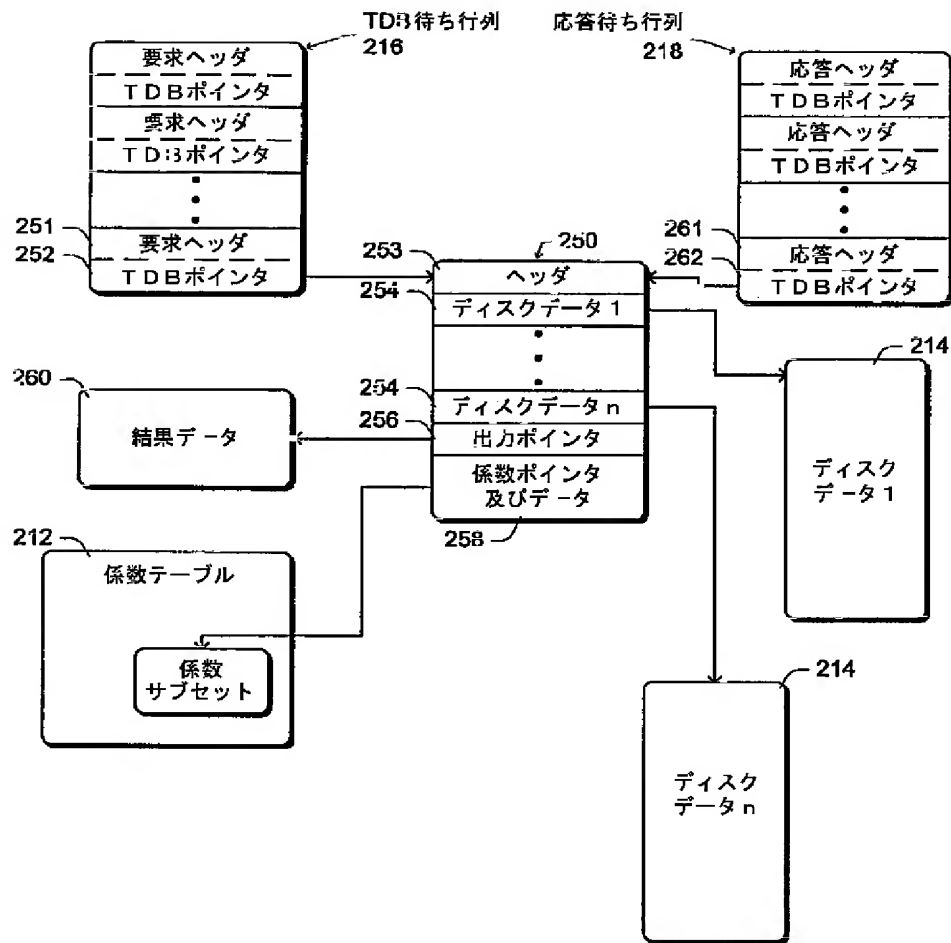
【図4】



【図5】



【図6】



マイクロプロセッサ

ハードウェア論理回路

